

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 10-083257

(43)Date of publication of application : 31.03.1998

(51)Int.Cl.

G06F 3/06
G06F 9/445
G06F 11/16

(21)Application number : 08-238013

(71)Applicant : MITSUBISHI ELECTRIC CORP

(22)Date of filing : 09.09.1996

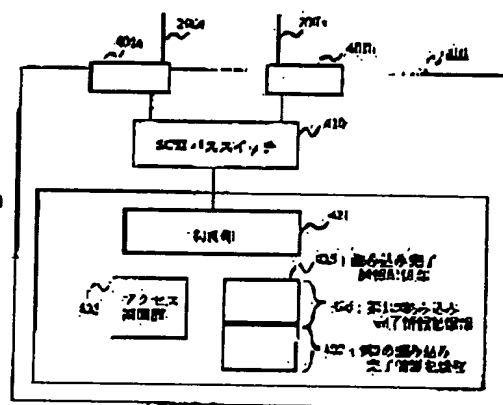
(72)Inventor : ITO KAZUHIKO
MIZUNO MASAHIRO
YAMAMOTO HITOSHI

(54) DATA STORAGE SYSTEM AND DATA STORAGE MANAGING METHOD

(57)Abstract:

PROBLEM TO BE SOLVED: To obtain a data storage system which can start two host systems cooperatively without a conflict in a data storage system shared by two host system.

SOLUTION: A decision part 421 discriminates which host system gains access. An access control part 423 exclusively control the access to a disk by the host system decided by the decision part 421 so that no conflict is caused. Further, a read completion information storage part 425 makes the two host systems know the completion of the system starting mutually.



LEGAL STATUS

[Date of request for examination] 30.10.1996

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 2830857

[Date of registration] 25.09.1998

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平10-83257

(43) 公開日 平成10年(1998) 3月31日

(51) Int.Cl. ⁶	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 3/06	3 0 1		G 0 6 F 3/06	3 0 1 C
9/445			11/16	3 1 0 Z
11/16	3 1 0		9/06	4 2 0 K

審査請求 有 請求項の数14 O L (全 13 頁)

(21) 出願番号 特願平8-238013

(22) 出願日 平成8年(1996) 9月9日

(71) 出願人 000008013

三菱電機株式会社

東京都千代田区丸の内二丁目2番3号

(72) 発明者 伊藤 一彦

東京都千代田区丸の内二丁目2番3号 三

菱電機株式会社内

(72) 発明者 水野 正博

東京都千代田区丸の内二丁目2番3号 三

菱電機株式会社内

(72) 発明者 山本 整

東京都千代田区丸の内二丁目2番3号 三

菱電機株式会社内

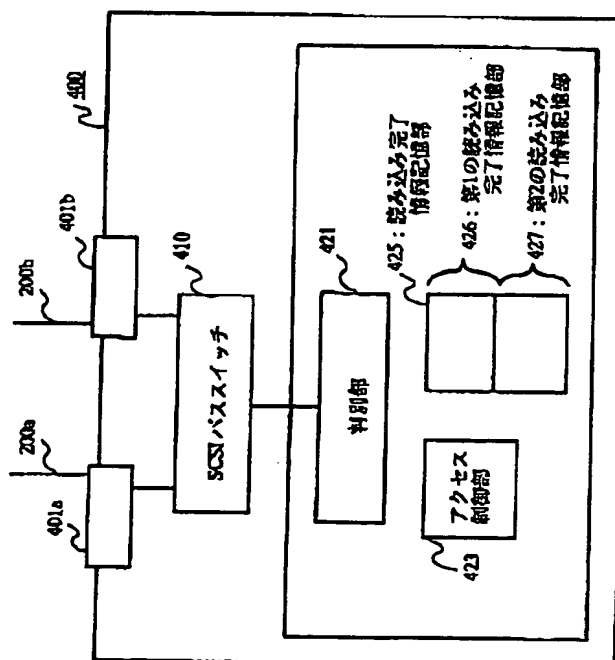
(74) 代理人 弁理士 宮田 金雄 (外3名)

(54) 【発明の名称】 データストレージシステム及びデータストレージ管理方法

(57) 【要約】

【課題】 2つのホストシステムで共有されるデータストレージシステムにおいて、2つのホストシステムに競合を発生させないように、協調してシステム起動を行えるデータストレージシステムを得る。

【解決手段】 判別部421がどちらのホストシステムからのアクセスかを判別する。アクセス制御部423は、判別部421により判別されたホストシステムのディスクへのアクセスを排他的に制御し、競合を発生させない。また、読み込み完了情報記憶部425により、2つのホストシステムは、お互いにシステム起動の完了を知る。



【特許請求の範囲】

【請求項1】 所定のオペレーティングシステムで動作する第1と第2のホストシステムに接続され、上記第1と第2のホストシステムからアクセスされるデータストレージシステムにおいて、

以下の要素を有するデータストレージシステム

(a) 所定のインタフェースを持つ第1と第2のバス、
(b) 上記第1のバスを介して上記第1のホストシステムに接続され、上記第2のバスを介して上記第2のホストシステムに接続されるとともに上記所定のオペレーティングシステムを記憶するストレージ、(c) 上記第1のバスを介して上記第1のホストシステムに接続され、上記第2のバスを介して上記第2のホストシステムに接続され、上記第1と第2のホストシステムの起動時に上記第1と第2のホストシステムによる上記ストレージに記憶されたオペレーティングシステムの読み込みを排他的に行わせるストレージマネージャ。

【請求項2】 上記ストレージマネージャは、上記第1と第2のホストシステムの起動時に、上記ストレージより先に上記第1と第2のホストシステムからアクセスされるように構成され、上記ストレージマネージャは、一方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込みを優先させて実行させ、他方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込みを待たせるとともに、

上記一方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込み完了後に、上記他方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込みを実行することにより上記第1と第2のホストシステムのオペレーティングシステムの読み込みを排他的に行わせることを特徴とする請求項1記載のデータストレージシステム。

【請求項3】 上記ストレージマネージャは、上記第1と第2のバスに接続され上記第1と第2のバスを切換えることにより上記第1と第2のバスのいずれかを選択的に接続すると共に、上記第1と第2のバスのどちらが接続されているかを示すバス接続情報を出力するバス切換機構と、

上記バス切換機構が出力するバス接続情報を参照し、上記第1と第2のホストシステムのどちらからアクセスされたかを判別する判別部と、

上記判別部の判別にに基づき、上記第1と第2のホストシステムの上記ストレージへのアクセスを排他制御し、上記ストレージに記憶されたオペレーティングシステムの読み込みを排他的に行うアクセス制御部を備えたことを特徴とする請求項2記載のデータストレージシステム。

【請求項4】 上記アクセス制御部は、上記第1と第2のバスのいずれか一方のバスをフリーズすることを特徴とする請求項3記載のデータストレージシステム。

【請求項5】 上記ストレージマネージャは、少なくとも上記第1と第2の一方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込みの完了を示す読み込み完了情報を記憶する読み込み完了情報記憶部を備え、上記アクセス制御部は、上記読み込み完了情報記憶部を参照して、上記一方のホストシステムのオペレーティングシステムの読み込み完了後に上記ストレージへのアクセスの排他制御を解除することを特徴とする請求項3に記載のデータストレージシステム。

【請求項6】 上記読み込み完了情報記憶部は、第1と第2の読み込み完了情報記憶部を備え、上記第1と第2の読み込み完了情報記憶部は、上記第1と第2のホストシステムのオペレーティングシステムの読み込みの完了を示す読み込み完了情報をそれぞれ個別に記憶するとともに、上記アクセス制御部は、第1と第2のホストシステムのいずれか先にアクセスがあったホストシステムのオペレーティングシステムの読み込みの完了まで、他方の後にアクセスのあったホストシステムのオペレーティングシステムの読み込みを待たせることを特徴とする請求項5に記載のデータストレージシステム。

【請求項7】 上記ストレージマネージャは、上記アクセス制御部により上記第1と第2のバスのうち一方のバスをフリーズすることにより、上記第1と第2のホストシステムのうちフリーズしていない他方のバスに接続されたホストシステムに上記ストレージに記憶されたオペレーティングシステムを読み込む優先権を与え、上記他方のバスに接続されたホストシステムは、オペレーティングシステムの読み込みの完了時に上記第1と第2の読み込み完了情報記憶部の一方に読み込み完了情報を記憶し、

上記アクセス制御部は、上記一方の読み込み完了情報記憶部を参照して、上記一方のバスのフリーズを解除することを特徴とする請求項6に記載のデータストレージシステム。

【請求項8】 上記ストレージは、ディスク装置であり、上記所定のインタフェースは、スモール・コンピュータ・システム・インタフェース（SCSI）であり、上記アクセス制御部は、スモール・コンピュータ・システム・インタフェース（SCSI）のディスク起動コマンドに対する応答を遅延することによりバスをフリーズすることを特徴とする請求項4記載のデータストレージシステム。

【請求項9】 所定のオペレーティングシステムの読み込みにより起動する第1と第2のホストシステムと上記オペレーティングシステムを記憶するデータストレージシステムとが接続され、上記データストレージシステムがストレージとストレージマネージャを有しているデータストレージシステムのデータストレージ管理方法において、

以下の工程を有するデータストレージ管理方法

(a) 上記第1と第2のホストシステムのいずれか一方のホストシステムが上記ストレージからオペレーティングシステムを読み込む工程、(b) 上記ストレージマネージャが他方のホストシステムの上記ストレージへのアクセスをロックするアクセスロック工程、(c) 上記一方のホストシステムによるオペレーティングシステムの読み込みの完了を検出し、アクセスロック工程による他方のホストシステムのアクセスのロックを解除するロック解除工程、(d) ロックを解除された他方のホストシステムが上記ストレージからオペレーティングシステムの読み込みを行う工程。

【請求項10】 上記アクセスロック工程は、(a) 他方のホストシステムの上記ストレージへのアクセス要求を検出するアクセス要求検出工程、(b) 上記アクセス要求検出工程により検出されたアクセス要求に対する応答を遅延する応答遅延工程、を有することを特徴とする請求項9記載のデータストレージ管理方法。

【請求項11】 上記データストレージ管理方法は、更に、上記一方のホストシステムのオペレーティングシステムの読み込み完了を示す読み込み完了情報を記憶する読み込み完了情報記憶工程を有し、上記ロック解除工程は、上記読み込み完了情報記憶工程により記憶された読み込み完了情報を参照して、他方のホストシステムのアクセスのロックを解除することを特徴とする請求項9記載のデータストレージ管理方法。

【請求項12】 第1と第2のホストシステムに接続され、上記第1と第2のホストシステムからアクセスされるデータストレージシステムにおいて、以下の要素を有するデータストレージシステム

(a) 所定のインタフェースを持つ第1と第2のバス、
(b) 上記第1のバスを介して上記第1のホストシステムに接続され、上記第2のバスを介して上記第2のホストシステムに接続されるとともにデータを記憶するストレージ、
(c) 上記第1と第2のホストシステムの上記ストレージに対するアクセスの可否を示すアクセス情報を予め設定するアクセス情報設定部、
(d) 上記アクセス情報設定部に設定されたアクセス情報を参照し、上記第1のホストシステムの上記ストレージに対するアクセスを制御する第1のアクセス制御部、
(e) 上記アクセス情報設定部に設定されたアクセス情報を参照し、上記第2のホストシステムの上記ストレージに対するアクセスを制御する第2のアクセス制御部。

【請求項13】 上記データストレージシステムは、更に上記アクセス情報を変更するアクセス情報変更部を備え、
上記第1と第2のアクセス制御部の少なくともいずれかは、一定の時間間隔で上記アクセス情報を参照することにより、上記第1と第2のホストシステムのうち少なくともいずれかのホストシステムの上記ストレージに対するアクセスを動的に制御することを特徴とする請求項1

2に記載のデータストレージシステム。

【請求項14】 上記データストレージシステムは、上記第1と第2のホストシステムの少なくともいずれか一方のホストシステムの障害を検出する障害検出部を備え、上記アクセス情報変更部は、上記障害検出部が障害を検出した一方のホストの上記ストレージへのアクセスを禁止すると共に、他方のホストの上記ストレージへのアクセスを許可するよう上記アクセス情報を変更することを特徴とする請求項13記載のデータストレージシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 この発明は、2つのホストシステムで共有されるデータストレージシステムに関するものである。また、2つのホストシステムで共有されるデータストレージ管理方法に関するものである。

【0002】

【従来の技術】

従来例1. 図13は、1つのディスクサブシステムを2つのホストシステムで共有するシステムの構成図である。図13において、100aは現用系(主系)のホストシステム、100bは予備系(従系)のホストシステムである。110aはホストシステム100aのディスク制御装置、110bはホストシステム100bのディスク制御装置である。ディスク300は、ホストシステム100a、100bで共有される複数のディスクである。200aは、ディスク300とホストシステム100aを接続するバス、200bは、ディスク300とホストシステム100bを接続するバスである。このような構成を採るシステムの一例として、デュプレクスシステム(Duplex System)がある。デュプレクスシステムは、スタンバイシステムとも呼ばれ、オンライン処理を行う現用系と予備系のプロセッサから構成される。予備系は、通常バッチ処理などの優先度の低い業務に使用される。現用系が障害になると予備系の業務を中止し、障害装置を予備に切り換えて処理を再開する。予備系を常時サービス開始可能な状態に待機させておく方式をホットスタンバイシステムと呼ぶ(電子情報通信ハンドブック、電子情報通信学会編、1988年3月30日発行による)。また、予備系をサービス開始可能な状態に待機させておかない方式もある。この方式をコールドスタンバイシステム(Cold Stand by system)と呼ぶ。コールドスタンバイシステムの場合には、現用系が障害になった時点で、OS(オペレーティングシステム)の読み込み等の起動操作を行い、サービス開始可能な状態とする。即ち、コールドスタンバイシステムの場合には、現用系ホストシステムの起動がまず行われ、現用系が障害になった時点で、予備系ホストシステムの起動が行われる。このため、それぞれのホストシステムの起動時には、ホストシ

システムとディスクサブシステムが1対1となるので通常のシステムと同様に起動動作が行われる。一方、ホットスタンドバイシステムの場合には、共有しているディスクサブシステムに対して主系・従系ともにアクセスを行うため、特に起動時に以下の問題が発生する。その問題とは、ディスクサブシステムに対して従系がアクセスしている間に主系のアクセスが受け付けられず、主系が動作を中断してハングアップ(Hung Up)してしまう可能性があることである。なぜならば、システム起動中は、ホストシステムがBIOS(Basic Input Output System)レベルで動作している為、OSの下で動作する本来のディスク制御用のドライバがロードされているときのようにリトライ動作等が十分に行われない為である。上記の理由から、デュプレクスシステム構成をとることが難しいという問題点があった。また、デュプレクスシステム構成でも従系をコールドスタンドバイにするなどの機能的な縮退が必要であるという問題点があった。また、デュプレクスシステムで従系をホットスタンドバイとする為には、主系、従系の間で協調してシステム起動を行うような方式をとる必要があった。

【0003】従来例2. 次に、上記従来例1と同一の構成のシステムにおける他の問題点について述べる。主系(または従系)に障害が生じた場合、両系による運転から従系(または主系)のみの運転に切り替えるシステムにおいて以下の問題が存在する。第1の問題は、障害が生じた主系(または従系)の運転をシャットダウンする場合に、障害が生じた主系(または従系)がデータストレージにもアクセスするため、同じデータストレージをアクセスしている従系(または主系)の運転に干渉する可能性があるという点である。第2の問題は、従系(または主系)のみの運転中(縮退運転中)、メンテナンスの為に主系(または従系)を立ち上げた場合に、従系(または主系)の運転に干渉する可能性があるという点である。このように、ディスクサブシステムを共有する一方の系の動作が他の系に影響を及ぼし、場合によっては、システムダウンに陥る可能性があるという問題点があった。

【0004】

【発明が解決しようとする課題】この発明は、上記のような問題点を解決するためになされたものであり、2つのホストシステムで共有されるデータストレージシステムにおいて、システム起動時に主系・従系の間で競合を発生させない様に協調してシステム起動を行えるデータストレージシステム及びその管理方法を得ることを目的としている。また、2つのホストシステムで共有されるデータストレージシステムにおいて、2つのホストシステムの切替時及び一方のホストシステムの保守時のデータストレージへの干渉を防ぐ機構を備えたデータストレージシステムを得ることを目的としている。

【0005】

【課題を解決するための手段】この発明に係るデータストレージシステムは、所定のオペレーティングシステムで動作する第1と第2のホストシステムに接続され、上記第1と第2のホストシステムからアクセスされるデータストレージシステムにおいて、以下の要素を有することを特徴とする。

(a) 所定のインタフェースを持つ第1と第2のバス、
(b) 上記第1のバスを介して上記第1のホストシステムに接続され、上記第2のバスを介して上記第2のホストシステムに接続されるとともに上記所定のオペレーティングシステムを記憶するストレージ、(c) 上記第1のバスを介して上記第1のホストシステムに接続され、上記第2のバスを介して上記第2のホストシステムに接続され、上記第1と第2のホストシステムの起動時に上記第1と第2のホストシステムによる上記ストレージに記憶されたオペレーティングシステムの読み込みを排他的に行わせるストレージマネージャ。

【0006】上記ストレージマネージャは、上記第1と第2のホストシステムの起動時に、上記ストレージより先に上記第1と第2のホストシステムからアクセスされるように構成され、上記ストレージマネージャは、一方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込みを優先させて実行させ、他方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込みを待たせるとともに、上記一方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込み完了後に、上記他方のホストシステムの上記ストレージに記憶されたオペレーティングシステムの読み込みを実行することにより上記第1と第2のホストシステムのオペレーティングシステムの読み込みを排他的に行わせることを特徴とする。

【0007】上記ストレージマネージャは、上記第1と第2のバスに接続され上記第1と第2のバスを切替えることにより上記第1と第2のバスのいずれかを選択的に接続すると共に、上記第1と第2のバスのどちらが接続されているかを示すバス接続情報を出力するバス切替機構と、上記バス切替機構が出力するバス接続情報を参照し、上記第1と第2のホストシステムのどちらからアクセスされたかを判別する判別部と、上記判別部の判別に基づき、上記第1と第2のホストシステムの上記ストレージへのアクセスを排他制御し、上記ストレージに記憶されたオペレーティングシステムの読み込みを排他的に行うアクセス制御部を備えたことを特徴とする。

【0008】上記アクセス制御部は、上記第1と第2のバスのいずれか一方のバスをフリーズすることを特徴とする。

【0009】上記ストレージマネージャは、少なくとも上記第1と第2の一方のホストシステムの上記ストレ

ジに記憶されたオペレーティングシステムの読み込みの完了を示す読み込み完了情報を記憶する読み込み完了情報記憶部を備え、上記アクセス制御部は、上記読み込み完了情報記憶部を参照して、上記一方のホストシステムのオペレーティングシステムの読み込み完了後に上記ストレージへのアクセスの排他制御を解除することを特徴とする。

【0010】上記読み込み完了情報記憶部は、第1と第2の読み込み完了情報記憶部を備え、上記第1と第2の読み込み完了情報記憶部は、上記第1と第2のホストシステムのオペレーティングシステムの読み込みの完了を示す読み込み完了情報をそれぞれ個別に記憶するとともに、上記アクセス制御部は、第1と第2のホストシステムのいずれか先にアクセスがあったホストシステムのオペレーティングシステムの読み込みの完了まで、他方の後にアクセスのあったホストシステムのオペレーティングシステムの読み込みを待たせることを特徴とする。

【0011】上記ストレージマネージャは、上記アクセス制御部により上記第1と第2のバスのうち一方のバスをフリーズすることにより、上記第1と第2のホストシステムのうちフリーズしていない他方のバスに接続されたホストシステムに上記ストレージに記憶されたオペレーティングシステムを読み込む優先権を与え、上記他方のバスに接続されたホストシステムは、オペレーティングシステムの読み込みの完了時に上記第1と第2の読み込み完了情報記憶部の一方に読み込み完了情報を記憶し、上記アクセス制御部は、上記一方の読み込み完了情報記憶部を参照して、上記一方のバスのフリーズを解除することを特徴とする。

【0012】上記ストレージは、ディスク装置であり、上記所定のインタフェースは、スモール・コンピュータ・システム・インタフェース（SCSI）であり、上記アクセス制御部は、スモール・コンピュータ・システム・インタフェース（SCSI）のディスク起動コマンドに対する応答を遅延することによりバスをフリーズすることを特徴とする。

【0013】この発明に係るデータストレージ管理方法は、所定のオペレーティングシステムの読み込みにより起動する第1と第2のホストシステムと上記オペレーティングシステムを記憶するデータストレージシステムとが接続され、上記データストレージシステムがストレージとストレージマネージャを有しているデータストレージシステムのデータストレージ管理方法において、以下の工程を有することを特徴とする。

（a）上記第1と第2のホストシステムのいずれか一方のホストシステムが上記ストレージからオペレーティングシステムを読み込む工程、（b）上記ストレージマネージャが他方のホストシステムの上記ストレージへのアクセスをロックするアクセスロック工程、（c）上記一方のホストシステムによるオペレーティングシステムの

読み込みの完了を検出し、アクセスロック工程による他方のホストシステムのアクセスのロックを解除するロック解除工程、（d）ロックを解除された他方のホストシステムが上記ストレージからオペレーティングシステムの読み込みを行う工程。

【0014】上記アクセスロック工程は、（a）他方のホストシステムの上記ストレージへのアクセス要求を検出するアクセス要求検出工程、（b）上記アクセス要求検出工程により検出されたアクセス要求に対する応答を遅延する応答遅延工程、を有することを特徴とする。

【0015】上記データストレージ管理方法は、更に、上記一方のホストシステムのオペレーティングシステムの読み込み完了を示す読み込み完了情報を記憶する読み込み完了情報記憶工程を有し、上記ロック解除工程は、上記読み込み完了情報記憶工程により記憶された読み込み完了情報を参照して、他方のホストシステムのアクセスのロックを解除することを特徴とする。

【0016】この発明に係るデータストレージシステムは、第1と第2のホストシステムに接続され、上記第1と第2のホストシステムからアクセスされるデータストレージシステムにおいて、以下の要素を有することを特徴とする。

（a）所定のインタフェースを持つ第1と第2のバス、

（b）上記第1のバスを介して上記第1のホストシステムに接続され、上記第2のバスを介して上記第2のホストシステムに接続されるとともにデータを記憶するストレージ、（c）上記第1と第2のホストシステムの上記ストレージに対するアクセスの可否を示すアクセス情報を予め設定するアクセス情報設定部、（d）上記アクセス情報設定部に設定されたアクセス情報を参照し、上記第1のホストシステムの上記ストレージに対するアクセスを制御する第1のアクセス制御部、（e）上記アクセス情報設定部に設定されたアクセス情報を参照し、上記第2のホストシステムの上記ストレージに対するアクセスを制御する第2のアクセス制御部。

【0017】上記データストレージシステムは、更に上記アクセス情報を変更するアクセス情報変更部を備え、上記第1と第2のアクセス制御部の少なくともいずれかは、一定の時間間隔で上記アクセス情報を参照することにより、上記第1と第2のホストシステムのうち少なくともいずれかのホストシステムの上記ストレージに対するアクセスを動的に制御することを特徴とする。

【0018】上記データストレージシステムは、上記第1と第2のホストシステムの少なくともいずれか一方のホストシステムの障害を検出する障害検出部を備え、上記アクセス情報変更部は、上記障害検出部が障害を検出した一方のホストの上記ストレージへのアクセスを禁止すると共に、他方のホストの上記ストレージへのアクセスを許可するよう上記アクセス情報を変更することを特徴とする。

【0019】

【発明の実施の形態】

実施の形態1. 以降の実施の形態では、ストレージがディスクであり、ストレージマネージャがディスクマネージャである場合を例にとって説明する。また、ディスクサブシステムがホストシステムにSCSI (Small Computer System Interface) バスで接続される場合について説明する。図1は、この発明のデータストレージシステムの構成図である。図1に示すように、この実施の形態の構成は、図13に示した従来の構成にディスクマネージャ400を加えたものである。また、ディスク300は、主系・従系両方のホストシステム100a、100bにそれぞれロードされるオペレーティングシステム310を記憶している。また、この実施の形態の構成は、主系・従系の2つのホストシステムにそれぞれ接続される2本のSCSIバス200a、200bを備えている。他の構成については、図13に示した従来の構成と同様である。また、ディスクサブシステムを共有する2つのホストシステムは、デュプレクスシステムを構成することを想定しているが、2つのホストシステムが同一動作を行うデュアルシステムでも構わない。ディスクマネージャ400は、共用ディスクサブシステム上に接続されるディスクサブシステム管理機構である。ディスクマネージャ400は、ディスク300と同様に主系・従系の両方のホストシステムからアクセス可能なデバイスである。また、ディスクマネージャ400は、ホストシステムからは、ディスクと同等に見える。以下、主系のホストシステムを単に主系ともいう。また、従系のホストシステムを単に従系ともいう。

【0020】図2は、ディスクマネージャ400の構成図である。ディスクマネージャ400は、主系・従系の2つのホストシステムにそれぞれ接続される2本のSCSIバス200a、200bに接続されている。SCSIバスは、それぞれ主系用ポート401a、従系用ポート401bを経て、SCSIバススイッチ410に接続される。このように、ディスクマネージャ400は、必要に応じて2つのSCSIバスを切り替えるデュアルポートディスク (Dual Port Disk) の構成を有する。SCSIバススイッチ410は、SCSIバスコントロール420に接続されている。また、ディスクマネージャ400は、MPU (Micro Processing Unit) 440と、不揮発性のメモリ430を有している。このように、ディスクマネージャ400は、SCSIバス上からは、ディスクとして認識され、ディスクと同等の動作を実行することが可能な構成となっている。SCSIバス上からアクセスされるデータは、ディスクマネージャ400上のメモリ430のSCSIブロックエリア432に記録される。メモリ430は、ディスクマネージャプログラムエリア438を

有している。ディスクマネージャプログラムエリア438は、MPU 440で実行されディスクマネージャ400を制御するディスクマネージャプログラムを記憶する。

【0021】図3は、この発明のデータストレージシステムの機能ブロック図である。図3に示す判別部421は、バス切換機構であるSCSIバススイッチ410からのポート番号を受け取り、受け取ったポート番号からどちらのバスからアクセスされたかを判別する。これにより、どちらのホストシステムからアクセスされたかが判別される。アクセス制御部423は、判別部421により判別されたホストシステムのディスクへのアクセスを制御する。読み込み完了情報記憶部425は、第1の読み込み完了情報記憶部426と、第2の読み込み完了情報記憶部427からなる。アクセス制御部423は、後述するように、この読み込み完了情報記憶部425を参照して、各ホストシステムのディスクへのアクセスを制御する。

【0022】次に、図4から図6を用いて、動作について具体的に説明する。この実施の形態では、システムに電源を投入する際、まず、主系のホストシステムから先に電源を投入し、次に、従系のシステムに電源を投入するものとする。また、主系、従系のシステム構成は、同一の構成であり、電源投入後のシステム起動にかかる時間もそれぞれ、ほぼ同一であることを想定している。また、ディスクマネージャ400は、ディスク300よりも先にホストシステム100a、100bからアクセスされる構成となっているものとする。図4は、主系のホストシステムの、電源投入からシステム起動完了までの流れ図である。図5は、主系と従系のホストシステムの電源投入からシステム起動完了までの、ディスクマネージャの動作の流れ図である。図6は、従系のホストシステムの、電源投入からシステム起動完了までの流れ図である。図4に示すように、まず主系がシステムのハードウェアへの電源投入 (S110) でシステム起動を開始する。次に、S120において、主系のホストシステムは、主系システムのハードウェア (H/W) をチェックする。チェック終了後、入出力制御装置上のBIOS (Basic Input Output System) によりディスクのリードライト (Read/Write) 動作を開始する (S130)。まず、主系のホストシステムは、ディスクサブシステムに対して、リセット (Reset) を発行する (S132)。次に、主系のホストシステムは、接続されているディスクを確認するコマンド (Target Select, Inquiry CMD, Test Unit Ready) を発行する (S134)。そして、主系のホストシステムは、ディスクのバス上をディスクマネージャから順番にサーチして、接続されているディスクそれぞれに対してディスクの識別子であるIDの所定の順 (昇順または降

順)にディスク起動コマンド(Start CMD)を発行していく。ディスクマネージャ400は、ホストシステムからはディスクとして認識される。また、ディスクマネージャ400は、ディスクよりも先に各ホストシステムからアクセスされる構成となっている。この為、主系のホストシステム100aは、まず最初にディスクマネージャ400にディスク起動コマンドを発行する(S136)。

【0023】一方、従系も、主系に続いて、システムのハードウェアへの電源投入でシステム起動を開始する

(図6、S310)。次に、S320において、従系のホストシステムは、従系システムのハードウェア(H/W)をチェックする。チェック終了後、入出力制御装置上のBIOS(Basic Input Output System)によりディスクのリードライト(Read/Write)動作を開始する(S330)。

まず、従系のホストシステムは、ディスクサブシステムに対して、リセット(Reset)を発行する(S332)。この時、すでに主系がディスクからOSの読み込みを実行中であれば、ディスクサブシステムに対して主系のホストシステムに接続されているバス200aが有効となる。主系のバス200aが有効である間、従系のバス200bからのアクセスは無視される。従って、ディスクサブシステムは、従系から発行されたリセットを無視する。このため、従系のリセット発行により、主系のOS読み込みが妨げられることはない。次に、従系のホストシステムは、接続されているディスクを確認するコマンドを発行する(S334)。そして、従系のホストシステムは、ディスクのバス上をディスクマネージャから順番にサーチして、接続されているディスクそれぞれに対してディスクの識別子であるIDの所定の順(昇順または降順)にディスク起動コマンドを発行していく。ディスクマネージャ400は、ホストシステムからはディスクとして認識される為、従系のホストシステムは、主系と同様に、まず最初にディスクマネージャ400にディスク起動コマンドを発行する(S336)。

【0024】ディスクマネージャ400は、主系・従系両方のホストシステムからディスク起動コマンドを受け取る。受け取ったディスク起動コマンドには、コマンドが主系用ポート、従系用ポートの内、どちらのポートから入ってきたかを示すポート番号がついている。ディスクマネージャ400の判別部421は、このポート番号を参照して、主系・従系のうちどちらからディスク起動コマンドが発行されたかを判別してアクセス制御部423へ出力する。アクセス制御部423は、ディスク起動コマンドが主系から発行された場合には、そのディスク起動コマンドに対して応答する。図5のS210に示すように、ディスクマネージャ400は、主系からのディスク起動コマンドを受け取ると、主系に対して、起動完了を報告する。また、S220に示すようにディスク起

動コマンドが従系から発行された場合には、ディスク起動コマンドを受け取った時点で従系のディスクのバス200bをフリーズ(FREEZE)する。バス200bがフリーズされると、従系はディスク起動コマンドの完了応答待ちになる(S337)。一般にディスク起動コマンドに対する応答のタイムアウトは設定されない為、従系は応答があるまで無制限に待つことになる。主系は、ディスクマネージャ400のアクセス制御部からディスク起動コマンドに対する起動完了を受け取ると、次に接続されているディスク300に対してディスク起動コマンドを発行する。その後、主系は、ディスクからOS(オペレーティングシステム)の読み込み(ロード)を行う(S138)。OSの読み込みが完了すると、主系のホストシステムのシステム起動が完了する。主系は、システムの起動完了後、ディスクマネージャ400に対して起動完了を報告する(S140)。起動完了の報告は、主系のホストシステムからディスクマネージャの管理するSCSIブロックエリアに完了コード(起動完了フラグ)を記録することで実行される。

【0025】図7は、ディスクマネージャのSCSIブロックエリアを示す図である。主系・従系ともにディスクマネージャ400上のSCSIブロックを共有する為、以下の方法で起動の完了報告を受け取ることができる。SCSIブロック0は、主系からは、リードライト(読み書き)可能なSCSIブロックであり、従系からはリード(読み込み)専用のSCSIブロックである。また、SCSIブロック1は従系からは、リードライト(読み書き)可能なSCSIブロックであり、主系からはリード(読み込み)専用のSCSIブロックである。主系・従系それぞれのホストシステムは、これらのSCSIブロックのうち書き込みする権利のあるSCSIブロックに起動完了フラグ433a、434aの書き込みを行う。他の系は、起動完了フラグの書き込みをチェックすることでシステム起動完了を検出できる。主系には、ディスクマネージャのSCSIブロックエリアの主系のホストシステム用に定められた所定の場所であるSCSIブロック0に対して、完了コード(例えば、数字の1)の書き込みを行う書き込み手段(プログラム)が用意されている。そして、その書き込み手段は、システムの起動完了後に自動的に起動されるよう予め設定されているものとする。その設定に従い、書き込み手段は、完了コードをディスクマネージャのSCSIブロックエリアの主系のホストシステム用に定められた所定の場所であるSCSIブロック0に記録する。主系のホストシステムは、その後、適当な時間間隔でSCSIブロック1をポーリングしながら、従系のシステム起動完了を待つ(S150)。

【0026】アクセス制御部423は、一定の時間間隔で、ディスクマネージャのSCSIブロックエリアの主系のホストシステム用に定められた所定の場所(SCS

Iブロック0)をチェックしている。アクセス制御部は、そのチェックにより、主系のシステム起動が完了したことを確認し(S230)、従系のバスのフリーズを解除し、ディスク起動コマンドに対する応答を返す(S240)。

【0027】従系は、ディスクマネージャ400のアクセス制御部423からディスク起動コマンドに対する起動完了を受け取ると、接続されているディスク300に対してディスク起動コマンドを発行する。その後、従系は、ディスクからOS(オペレーティングシステム)の読み込み(ロード)を行う(S338)。OSの読み込みが完了すると、従系のホストシステムのシステム起動が完了する。従系は、システムの起動完了後、ディスクマネージャ400に対して起動完了を報告する(S340)。従系の起動完了の報告も主系と同様に、従系のホストシステムからディスクマネージャの管理するSCSIブロックエリアに完了コードを記録することで実行される。主系と同様に、従系には、ディスクマネージャのSCSIブロックエリアの従系のホストシステム用に定められた所定の場所(SCSIブロック1)に対して、完了コード(例えば、数字の1)の書き込みを行う書き込み手段が用意されている。そして、その書き込み手段は、主系の書き込み手段と同様に、システムの起動完了後に自動的に起動されるよう予め設定されているものとする。その設定に従い、書き込み手段は、完了コードをディスクマネージャのSCSIブロックエリアの従系のホストシステム用に定められた所定の場所(SCSIブロック1)に記録する。SCSIブロック1をポーリングしている主系のホストシステムは、完了コードの書き込みによりシステム起動完了を知る。その後、主系と従系は、ディスクサブシステムを共有して、デュプレクスシステムの運用を開始する。

【0028】なお、ここでは、主系から発行されたディスク起動コマンドが先にディスクマネージャに受け取られる場合について、図5に示している。これは、前述したように、主系と従系のシステムが同一構成であり、主系のホストシステムの電源投入が先に行われるため、主系のディスク起動コマンドが先に発行されることを想定しているためである。だが、どちらのホストシステムからコマンドが発行されたかは、前述したように判別部により判別されるため、仮に従系のディスク起動コマンドが先に発行されても構わない。即ち、図5のS210とS220は順不同で構わない。従系のディスク起動コマンドが先に発行された場合には、従系のバスが主系のディスク起動コマンドの発行に先立ちフリーズされる。その後、主系からディスク起動コマンドが発行されると、主系に対して、起動完了が報告される。その後のOSの読み込み動作は、前述した場合と同様である。

【0029】次に、ディスクマネージャの動作について説明する。図8は、システム起動時のディスクマネー

ジャの動作を示す流れ図である。ディスクマネージャ400は、通常は主系・従系の動作に対する応答をSCSIバスを通じて行う。システム起動時には図8に示すシーケンスで動作を実行する。まず、S500において、ディスクマネージャ400の初期化が行われる。次に、S510において主系・従系両方のシステム起動が完了しているかをチェックする。完了していれば、S520において、通常動作に移行する。完了していなければ、S530において、主系のシステム起動が完了しているかどうかをチェックする。主系のシステム起動が完了していなければ、S540において、主系のコマンドに対する応答を優先して行う。このとき、従系からのディスク起動コマンドが発行されたら、アクセス制御部により、従系のSCSIバスをフリーズする。主系のシステム起動が完了していれば、S550において、起動コマンドをフリーズしているかどうか判定する。フリーズしていれば、S560において、起動コマンドに対する応答を実行する。フリーズしていなければ、従系は、コマンドの応答待ちの状態ではないので、S560の処理はスキップする。次に、S570において、従系のコマンドに対する応答を優先して行う。

【0030】このように、アクセス制御部は、システム起動時の主系と従系のディスクアクセスを制御し、OSの読み込みを排他的に行わせる。それにより、OSの読み込みで主系と従系の競合が発生しなくなる。競合が発生しないのでBIOSによるOSの読み込み時にエラーの発生を防ぐことができる。

【0031】実施の形態2. この実施の形態では、2重系のディスクに対する競合を制御するシステム管理機構を備えたデータストレージシステムについて説明する。図9は、この実施の形態のデータストレージシステムの構成図である。ディスクマネージャ400は、SCSI等のディスクインタフェースを有するディスクサブシステム管理機構(ハードウェア)である。ディスクマネージャ400は、主系と従系のホストシステムの各ディスクへのアクセスの許可(O)/禁止(X)をアクセス情報として登録するディスク管理テーブル1400を有している。アクセス制御部120a、120bはディスク管理テーブル1400を参照してホストシステム100a、100bのディスクへのアクセスを制御する。アクセス情報変更部122a、122bは、ディスク管理テーブル1400に登録されているアクセス情報を必要に応じて変更する。また、障害検出部124a、124bは、ホストシステム100a、100bに発生した障害を検出する。図10に、ディスク管理テーブルの一例を示す。図10に示すディスク管理テーブルでは、主系のホストシステムがディスク#0にアクセスを許可され(1410)、従系のホストシステムがディスク#1にアクセスを許可されている(1420)。ここでは、アクセス情報として、アクセス許可とアクセス禁止を

“○”と“×”で示しているが、他の値でも構わない。主系のディスク制御装置がディスク#0、従系のディスク制御装置がディスク#1を占有して主系と従系が動作する。このアクセスの可否を示すアクセス情報は、システム構築時にディスクマネージャのディスク管理テーブルへセットされる。ハードウェアの構成からはどちらのディスク制御装置からもディスク#0、#1共アクセス可能であるが、アクセスを許可されたディスク以外はアクセスしないように、アクセス制御部が制御を行なう。

【0032】アクセス制御部120a、120bの動作について説明する。ディスク管理テーブル1400には、システムの起動前に予めアクセス情報が設定されているものとする。ディスクマネージャ400は、SCSI等のディスクインタフェースを有する。アクセス制御部120a、120bはこのディスクインタフェースを介してディスク管理テーブル1400にアクセスすることができる。また、ディスクマネージャ400は、デュアルポート機構を持ち、2つの系からアクセスが可能である。アクセス制御部120a、120bは、ディスク管理テーブルを参照して、その系に対するディスクのアクセス許可または禁止を判別し、アクセスが禁止されているディスクへのアクセスを行なわない様に制御する。このように、ディスク管理テーブルを用いることで、物理的には、アクセス可能な複数のディスクを、ホストシステム毎に占有して使用することができる。

【0033】ディスク管理テーブルのアクセス情報は、アクセス情報変更部122a、122bにより必要に応じて随時変更される。一方、アクセス制御部120a、120bは、一定間隔でディスクマネージャ上のディスク管理テーブルを参照して、アクセス情報をチェックする。それにより、一旦システム起動した後でも、ディスクへのアクセスの制御（許可または禁止）をダイナミックに切り替えることができる。ディスク管理テーブルのアクセス情報の変更について具体的に説明する。システムに障害が生じた場合、ディスクアクセスを行わずに安全（他系に影響をおよぼさないよう）にシャットダウン動作を行なう為に、以下のようにディスク管理テーブルのアクセス情報を書き換える。ここでいうシステムに生じた障害とは、例えば、ローカルエリアネットワーク（LAN: Local Area Network）のネットワーク用制御カード（図示せず）に発生したトラブルなど、ホストシステムのCPUは正常に動作し、また、ホストシステムのアクセス制御部や、アクセス情報変更部及び障害検出部は動作可能であるような障害であるものとする。ここでは、障害が主系のホストシステム100aに発生した場合について説明する。障害検出部124aは、ホストシステム100aに発生した障害を検出する。障害検出部124aは障害が発生したことをアクセス情報変更部に知らせる。アクセス情報変更部は、主系のホストシステム100aからのアクセスが許

可になっているディスク#0のアクセス情報をアクセス禁止に書き換えて、主系のホストシステム100aのディスク#0へのアクセスを禁止すると共に、従系のホストシステム100bからのディスク#0へのアクセスを許可するように、ディスク管理テーブル1400のディスク#0のアクセス情報を書き換える。図11に、アクセス情報変更後のディスク管理テーブルを示す。図11の、1430に示すように主系からのディスク#0に対してはアクセス禁止となり、1440に示すように従系からディスク#0に対してはアクセス許可となっている。従系のアクセス制御部120bがディスク管理テーブルをポーリングして、ディスク#0が“アクセス禁止”から“アクセス許可”に変更になった事を知る。主系は、システムシャットダウン時、及び、診断用に再起動した場合も含めて、ディスク管理テーブルを参照する事によって、ディスクへのアクセスが禁止されている事を知り、ディスクへアクセスすることはない。従って、主系が従系の動作に干渉する事はない。

【0034】次に、主系のホストシステム100aがダウンし、主系のアクセス制御部120a、アクセス情報変更部122a、障害検出部124aが動作しない場合について説明する。この場合の、システムダウンの検出は、従系の障害検出部124bが行う。図12は、この一方の系の障害検出部が他方の系の障害を検出する場合のシステム構成図である。図12において、1100は、システムが生きているというパルスであるハートビートを主系と従系で送受信するハートビート信号線である。ハートビートは、システムが生きている間、一定の時間間隔で発生するパルスである。このハートビートを監視することにより、一定の時間間隔でハートビートが確認できれば、システムが正常稼働していることが判り、ハートビートが確認できなければ、なんらかの障害が発生したことを検出できる。障害検出部124bは、ハートビート信号線1100からハートビートを入力し監視することにより、主系のホストシステムのダウンを検出する。システムダウンを検出すると、障害検出部124bは、アクセス情報変更部122bに主系の障害発生を通知する。アクセス情報変更部122bは、主系のアクセスしていたディスク#0のアクセス情報の変更を行う。アクセス情報の変更については、すでに述べた変更と同様である。

【0035】また、アクセス情報変更部122a、122bが、ディスク管理テーブルのアクセス情報を内容に変更がなくとも一定の時間間隔で更新することにより、障害の発生を検出してもよい。アクセス情報変更部122a、122bが、アクセス情報を更新する際、更新のタイムスタンプ（日時の記録）も更新される。障害検出部124a、124bは、互いに相手のホストシステムのアクセス情報変更部が更新したタイムスタンプをチェックし、タイムスタンプが更新されていないとき、相手

のホストシステムに障害が発生したと判断する。障害発生後の処理は、すでに述べた処理と同様である。

【0036】以上のように、この実施の形態では、各ホストシステム毎にどのディスクにアクセスするかをアクセス情報としてディスク管理テーブルに設定し、アクセス制御部が、ディスク管理テーブルを参照して、各ホストシステムのディスクへのアクセスを制御するデータストレージシステムについて説明した。アクセス情報変更部は、障害検出部が障害を検出すると、障害を検出したホストシステムのディスクへのアクセスを禁止するようにアクセス情報を変更する。これにより、障害が発生したホストシステムのシャットダウンが、他のホストシステムに影響を与えることを回避する。また、アクセス情報変更部は、障害が検出されたホストシステムがアクセスしていたディスクを他のホストシステムからアクセスできるようにアクセス情報を変更する。これにより、ディスクが利用できない状態になることを回避できる。

【0037】

【発明の効果】この発明によれば、アクセス制御部が一方の系にOSをロードする優先権を与え、他の系からのアクセスをフリーズするので、読み込みの競合が発生せず、正常に読み込みを行うことができる。

【0038】また、この発明によれば、アクセス制御部が、ディスク管理テーブルを参照して、ホストシステムのディスクへのアクセスを制御する。障害発生時に、アクセス情報を変更し、障害の発生した系からディスクアクセスしないように制御することにより、正常な系に影響を及ぼすことが防止できる。

【図面の簡単な説明】

【図1】 この発明のデータストレージシステムの構成図である。

【図2】 この発明のデータストレージシステムのディスクマネージャ400の構成図である。

【図3】 この発明のデータストレージシステムの機能ブロック図である。

【図4】 この発明のデータストレージシステムの主系のホストシステムの電源投入からシステム起動完了まで

の流れ図である。

【図5】 この発明のデータストレージシステムの主系と従系のホストシステムの電源投入からシステム起動完了までのディスクマネージャの動作の流れ図である。

【図6】 この発明のデータストレージシステムの従系のホストシステムの電源投入からシステム起動完了までの流れ図である。

【図7】 この発明のデータストレージシステムのディスクマネージャのSCSIブロックエリアを示す図である。

【図8】 この発明のデータストレージシステムのシステム起動時のディスクマネージャの動作を示す流れ図である。

【図9】 この発明のデータストレージシステムの構成図である。

【図10】 この発明のデータストレージシステムのディスク管理テーブルの一例を示す図である。

【図11】 この発明のデータストレージシステムのアクセス情報変更後のディスク管理テーブルを示す図である。

【図12】 この発明のデータストレージシステムの他のシステム構成図である。

【図13】 従来の、1つのディスクサブシステムを2つのホストシステムで共有するシステムの構成図である。

【符号の説明】

100a、100b ホストシステム、200a、200b バス、300 ディスク、400 ディスクマネージャ、410 SCSIバススイッチ、420 SCSIバスコントロール、421 判別部、423 アクセス制御部、425 読み込み完了情報記憶部、426 第1の読み込み完了情報記憶部、427 第2の読み込み完了情報記憶部、430 メモリ、432 SCSIブロックエリア、438 ディスクマネージャプログラムエリア、440 MPU、1400 ディスク管理テーブル。

【図10】

ディスク管理テーブル			1420
	主系	従系	
ディスク#0	○	×	1420
ディスク#1	×	○	

○: アクセス許可
×: アクセス禁止

システム起動時のディスク管理テーブル

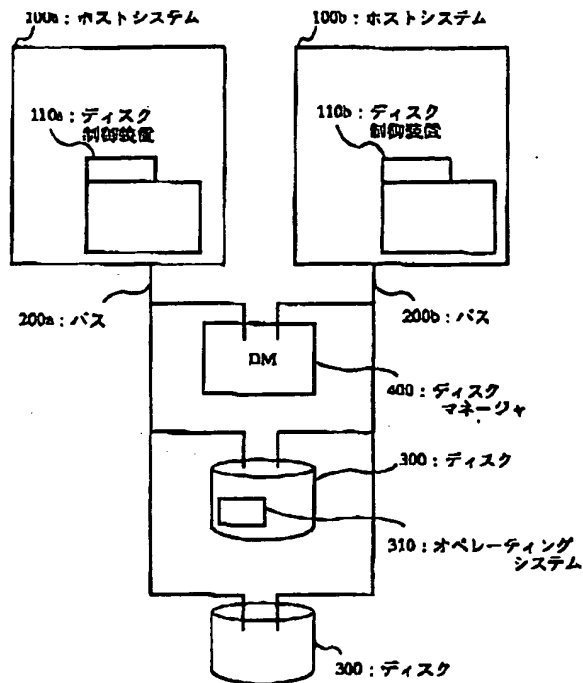
【図11】

ディスク管理テーブル			1440
	主系	従系	
ディスク#0	×	○	1440
ディスク#1	×	○	

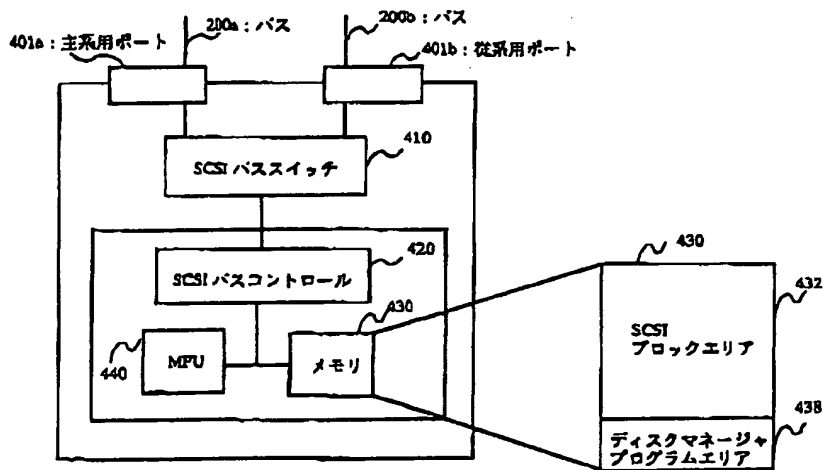
○: アクセス許可
×: アクセス禁止

主系故障時のディスク管理テーブル

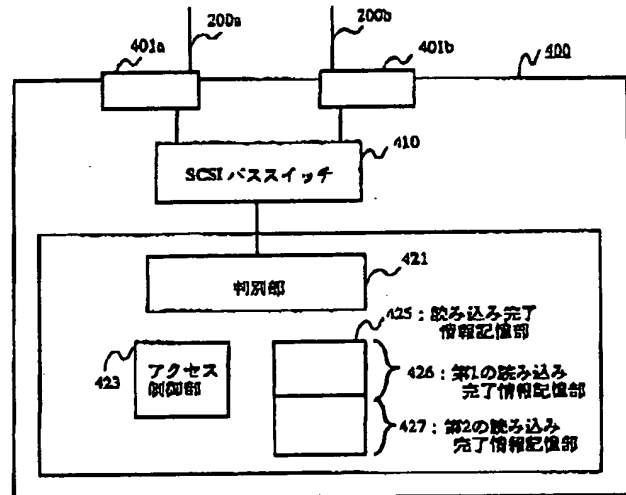
【図1】



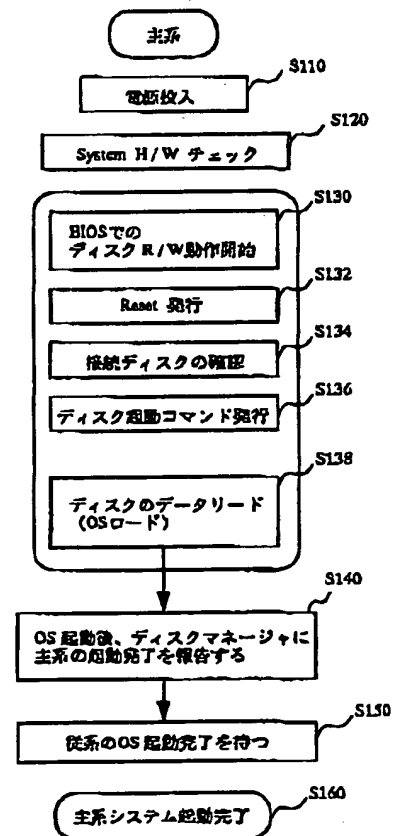
【図2】



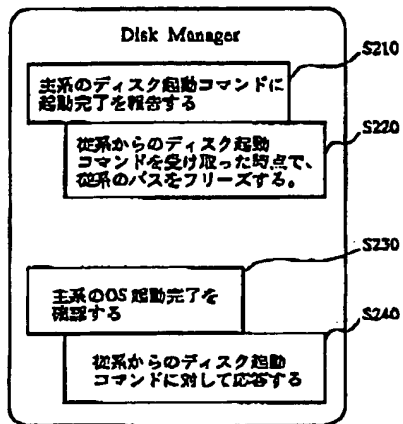
【図3】



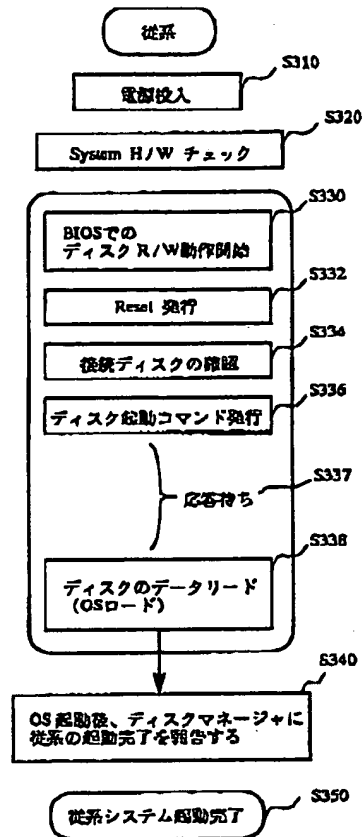
【図4】



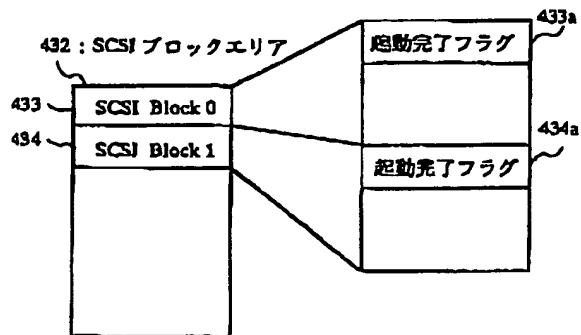
【図5】



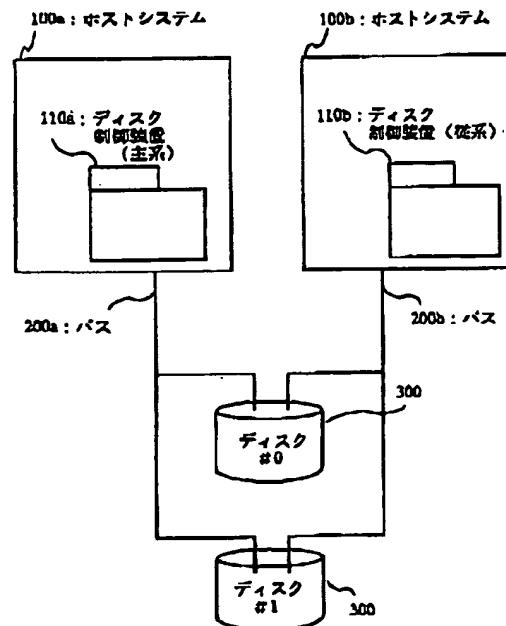
【図6】



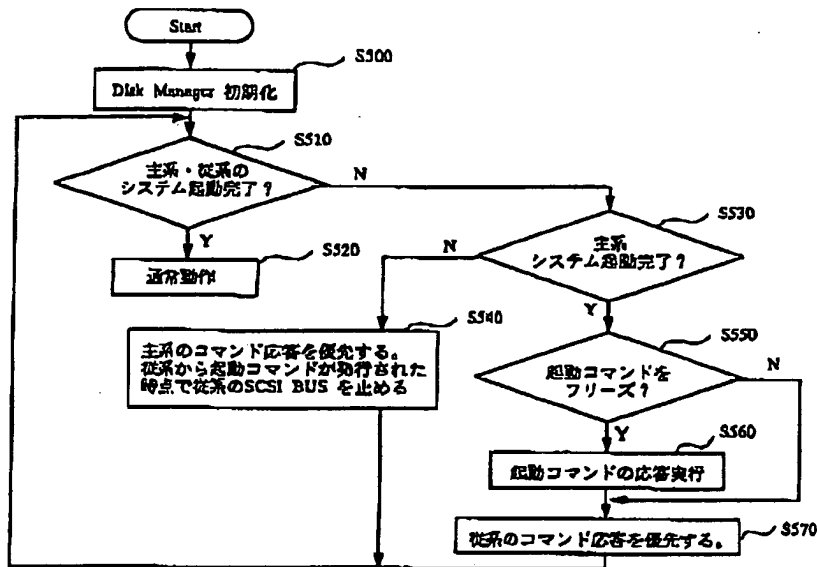
【図7】



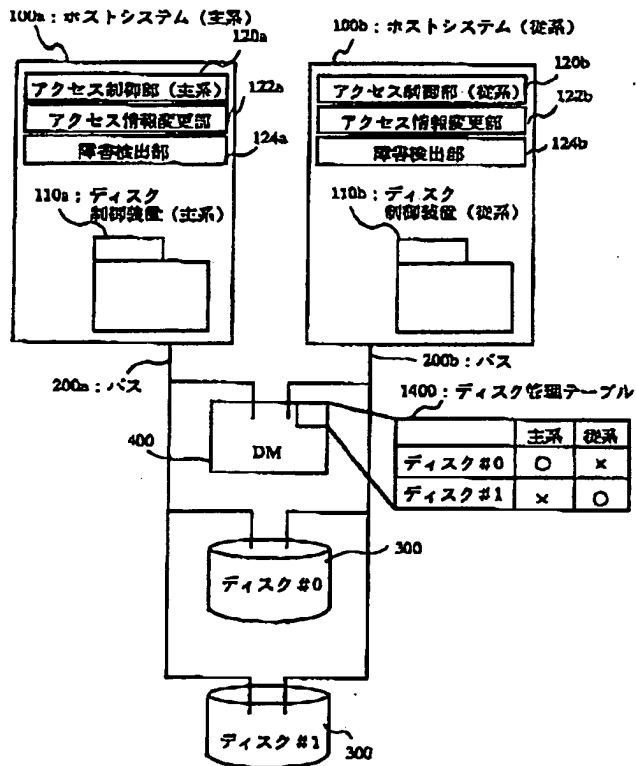
【図13】



【図8】



【図9】



【図12】

